

JCAHPCの次期スーパーコンピュータ Oakforest-PACS

東京大学 情報基盤センター JCAHPC 施設長

中村宏







最先端共同HPC基盤施設 JCAHPC

- Joint Center for Advanced High Performance Computing (http://jcahpc.jp)
- ・平成25年3月、筑波大学と東京大学は「計算科学・工学 及びその推進のための計算機科学・工学の発展に資す るための連携・協力推進に関する協定」を締結
- ・本協定の下、筑波大学計算科学研究センターと東京大学情報基盤センターが JCAHPC を設置
 - ・東京大学柏キャンパスの東京大学情報基盤センター内に、両機関の教職員が中心となって設計するスーパーコンピュータシステムを設置し、最先端の大規模高性能計算基盤を構築・運営するための組織







Oakforest-PACS in JCAHPC

- ・ 筑波大学と東京大学の間の密な連携・協力
- ・ 仕様を統一、計算資源として1つのシステム
- 2大学が調達と運用に関して責任を持つ
 - 国内初の試み
 - 日本で最大規模のシステムを実現





平成25年のプレスリリース

平成 25 年 7 月 22 日

国立大学法人 東京大学情報基盤センター 国立大学法人 筑波大学計算科学研究センター 最先端共同 HPC 基盤施設

最先端共同 HPC 基盤施設の活動を開始

筑波大学と東京大学によるスーパーコンピュータ共同開発、共同運営・管理





HPCI: High Performance Computing Infrastructure 日本全体におけるスパコンインフラ

今後のHPCI 計画推進の在り方について(H26/3)より

我が国の次期スパコン開発の方向性

< 我が国の計算科学技術インフラのイメージ>

http://www.mext.go.jp/b_menu/shingi/chousa/shinkou/028/gaiyou/1348991.htm

2020年頃までのエクサス

フラッグシップシステム
世界トップレベルの能力を有し幅広い分野をカバーするシステム

HPCIを通じて 国全体のインフ ラとして運用

フラッグシップを支える 特徴ある複数のシステム

9大学情報基盤センターのシステム 附置研・共同利用機関のシステム 独立行政法人のシステム

その他大学等のシステム

<u>ケールコンピューティン</u> <u>グ</u>の実現を目指す

リーディング マシン

能力を最大限に活かす アプリ開発も戦略的に推進







9大学情報基盤センター運用&整備計画(2016年5月時点)

Fiscal Year	2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025
Hokkaido III	HITACHI SR16000/M1 (172TF, 22TB) 3.2 PF (UCC + CFL/M) 0.96MW 30 PF (UCC + CFL-M) 2MW 30
Tohoku	NEC SX-9 (60TF) SX-ACE(707TF,160TB, 655TB/s) LX406e(31TF), Storage(4PB), 3D Vis, 2MW ~30PF, ~30PB/s Mem BW (CFL-D/CFL-M) ~3MW
Tsukuba	HA-PACS (1166 TF) PACS-X 10PF (TPF) 2MW COMA (PACS-IX) (1001 TF) Post T2K: Oakforest-PACS 25 PF 100+ PF 4.5MW
Tokyo	Column C
Tokyo Tech.	TSUBAME 2.5 (5.7 PF, 110+ TB, 1160 TB/s), 1.4MW TSUBAME 2.5 延長運転 3~4 PF TSUBAME 3.0 (20 PF, 4~6PB/s) 2.0MW TSUBAME 4.0 (100+ PF, (3.5アップグレード可能なら2018に40PF) >10 PB/s, ~2.0MW)
Nagoya -	FX10(90TF)
Kyoto	Cray: XE6 + GB8K + Cray XC40(5.5PF) + CS400(1.0PF) 50-100+ PF Cray XC30 (584TF) 1.3 MW (FAC/TPF + UCC) 1.8-2.4 MW
Osaka	NEC SX-ACE NEC Express5800 3.2PB/s,10~20Pflop/s, 1.0-1.5MW (CFL-M) 25.6 PB/s, 50-100Pflop/s,1.5-2.0MW 0.7-1PF (UCC)
Kyushu 📆	HA8000 (712TF, 242 TB) SR16000 (8.2TF, 6TB) 2.0MW FX10 (272.4TF, 36 TB) CX400 (966.2TF, 183TB) 15-20 PF (UCC/TPF) 2.6MW FX10 (90.8TFLOPS)

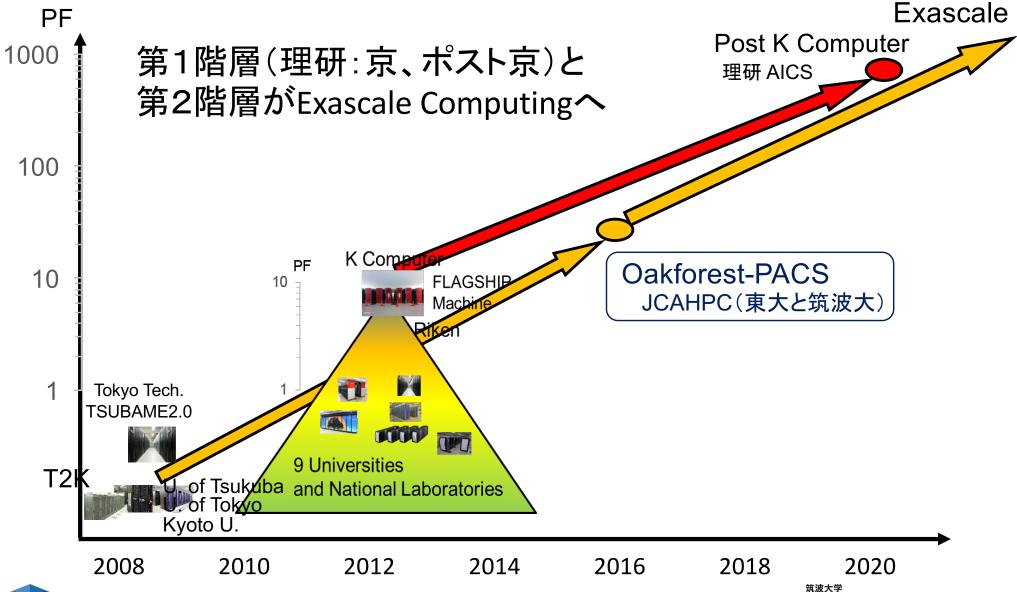






フラグシップとの両輪として

Future





JCAHPC: 共同調達への道のり

- 2013活動開始
 - 第1期(2013/4-2016/3):施設長:佐藤三久(筑波大学)、副施設長:石川裕(東京大学)
 - 第2期(2015/4-):施設長:中村宏(東京大学)、副施設長:梅村雅之(筑波大学)
- ・共同調達・運用へ向けて
 - 2013/7: RFI(request for information)共同調達は既定路線ではなかった→1システムとして調達へ
 - ・ 複数大学による初めての「1システム」共同調達へ
- どうして共同調達ができたのか?共同調達は大変・・
 - 目標を共有できる、ことに尽きる

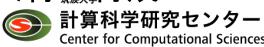






2センターのミッション

- ・ 筑波大学計算科学研究センターのミッション:
 - 計算機科学と計算科学の協働:学際的な高性能計算機開発
 → PACSシリーズの開発: CP-PACS@1996 TOP1
 - 先端学際科学共同研究拠点: 最先端の計算科学研究推進
 - これからの計算科学に必要な学際性を持つ人材を育成
- ・東京大学情報基盤センターのミッション:
 - 学際大規模情報基盤共同利用・共同研究拠点(8大学の情報基盤センター群からなるネットワーク型)の中核拠点: 大規模情報基盤を活用し学際研究を発展
 - HPCI資源提供機関:最先端スパコンの共同設計開発及び 運用、Capability資源および共用ストレージ資源の提供
 - 人材育成:計算科学の新機軸を創造できる人材の育成

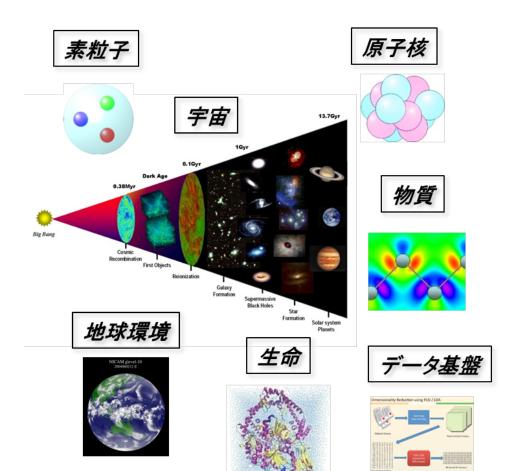


筑波大学計算科学研究センター^{♀♀」CAHPC}

計算科学と計算機科学の協働(コデザイン)

先端計算科学 推進室 次世代計算機システム 開発室

PACSシリーズの開発







筑波大学計算科学研究センター^{♀♀」CAHPC}

1992年4月 計算物理学研究センター設置(10年計画)。

1996年9月 CP-PACS(2048PU)完成, TOP500で世界第1位。

2004年4月 改組拡充し、計算科学研究センターを設置。

2007年4月 融合型宇宙シミュレータFIRST完成。

2008年6月 T2K-Tsukubaオープンスーパーコンピュータ運用開始。

2010年4月 共同利用・共同研究拠点「先端学際計算科学共同研究拠点」認定。

東京大学との協定に基づき「最先端共同HPC基盤施設」を設置。 2013年3月

科学者と計算機工学者の協力による。 application –drivenな超高速計算機の 開発・製作=学際計算科学

世界的に見てもユニーク

高い計算パワーの集中による計算科学 の最重点課題・最先端課題の研究

1978 1st PACS-9

1980 2nd PAXS-32



1989 5th QCDPAX







2012 8th HA-PACS



2014

9th COMA



8		
Year	System	Performance
1978	PACS-9 (PACS I)	7 KFLOPS
1980	PACS-32 (PACS II)	500 KFLOPS
1983	PAX-128 (PACS III)	4 MFLOPS
1004	- · - · - · · · · · · · · · · · · · · ·	A MEL ODG

	/	
1980	PACS-32 (PACS II)	500 KFLOPS
1983	PAX-128 (PACS III)	4 MFLOPS
1984	PAX-32J (PACS IV)	3 MFLOPS
1989	QCDPAX (PACS V)	14 GFLOPS
1996	CP-PACS (PACS VI)	614 GFLOPS
2006	PACS-CS (PACS VII)	14.3 TFLOPS
2012	HA-PACS (PACS VIII)	1.166 PFLOPS
2014	COMA (PACS IX)	1.001 PFLOPS

2007 **FIRST** (Hybrid Simulator)



36TFLOPS **Host 3TFOPS** Accelerator 33TFLOPS

筑波大学



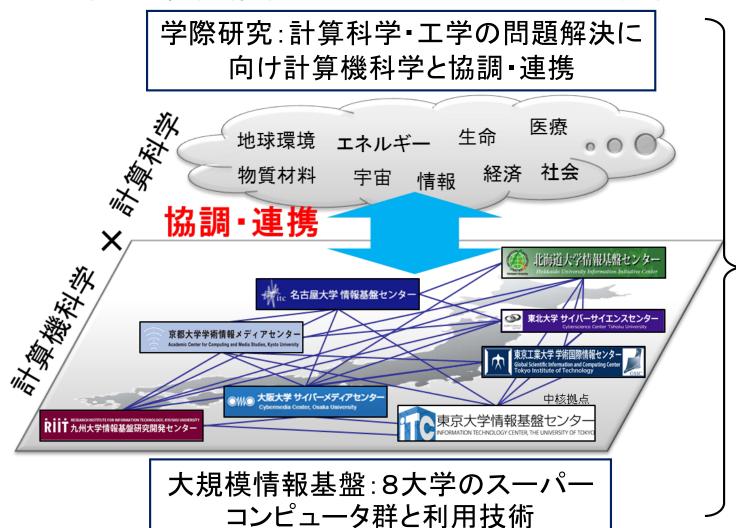




東京大学情報基盤センター

Research Center for Extreme Scale Computing and Data

学際大規模情報基盤共同利用・共同研究拠点の中核拠点



- ・解決や解明が困 難と考えられてい た課題の解決へ
- ・学術基盤としての 大規模情報基盤 の活用による 研究コミュニティへ の貢献
- 多様で大規模な計 算資源
- ・公募型の学際的 共同研究(萌芽段 階を含む)を遂行







東京大学情報基盤センター

- HPCI資源提供機関として
 - 機関連携による最先端スパコンの共同設計開発及び運用、Capability資源および共用ストレージ資源の提供
 - Data Intensive Applicationに対応したシステムの整備
- 人材育成機関として
 - 計算科学の新機軸を創造できる人材の育成
 - 学内各部局, 利用者, 共同利用・共同研究拠点との連携
- 合計約2,000人のユーザー(学外が半分)
 - 大学(研究,教育),研究機関,企業
- → 大規模シミュレーション, 特に連成解析
 - 全球規模大気海洋カップリング
 - ppOpen-HPC, ppOpen-MATH/MP
 - 地震シミュレーション
 - 地震発生+破壊伝播+強震動
 - 地盤強震動+都市・建造物振動
 - 流体・構造シミュレーション



Total Peak performance: 1.13 PFLOPS

Total number of nodes: 4800 Total memory: 150 TB

Peak performance / node: 236.5 GFLOPS

Main memory per node: 32 GB

Disk capacity: 1.1 PB + 2.1 PB

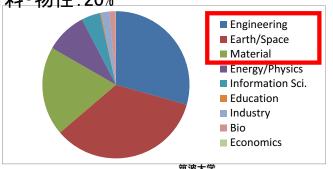
SPARC64 Ixfx 1.84GHz



利用の多い分野(2015年度)

- 工学(流体・構造・電磁気等):30%
- 地球宇宙科学(大気海洋・地震等):35%

- 材料・物性:20%



JCAHPC共同調達のポリシー ~2センターで共有したこと~



- T2Kの精神に基づき、オープンな最先端技術を導入
 - T2K: 2008年に始まったTsukuba, Tokyo, Kyoto の3大学でのオープンスパコンアライアンス、3機関の研究者が仕様策定に貢献、システムへの要求事項を共通化
- ・システムの基本仕様
 - ・ 超並列PCクラスタ
 - HPC用の最先端プロセッサ、アクセラレータは不採用
 - → 広範囲なユーザとアプリケーションのため
 - →ピーク性能追求より、これまでのコードの継承を優先
 - 使いやすい高効率相互結合網
 - 大規模共用ファイルシステム



Oakforest-PACS

- スケールメリットを活かす
 - 超大規模な単一ジョブ実行も可能とする





設置予定場所:東京大学柏キャンパス

Google マップ

https://www.google.com/maps/@?dg=dbrw&newdg=1





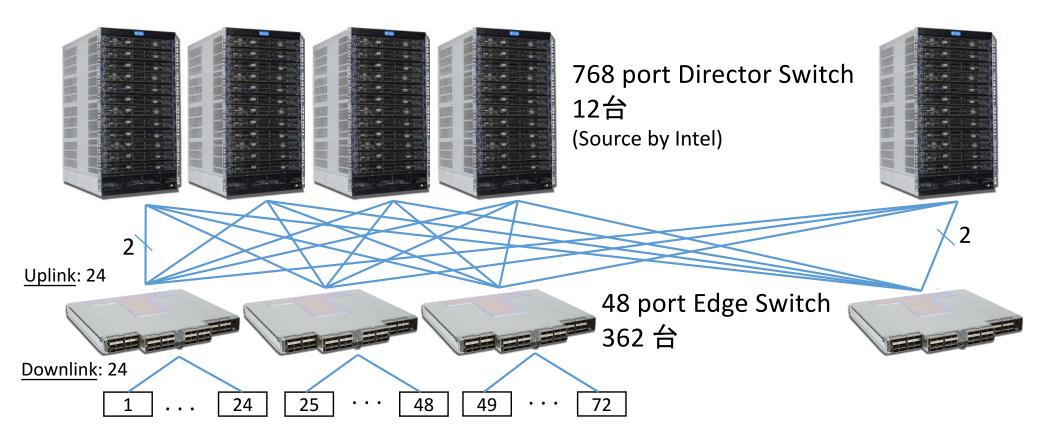
Oakforest-PACS の仕様

総ピーク演算性能			25 PFLOPS
ノード数			8,208
計算 ノード			富士通 PRIMERGY CX600 M1 (2U) + CX1640 M1 x 8node
	プロセッサ		Intel® Xeon Phi™ 7250 (開発コード: Knights Landing) 68 コア、1.4 GHz
	メモリ	高バンド幅	16 GB, MCDRAM, 実効 490 GB/sec
		低バンド幅	96 GB, DDR4-2400, ピーク 115.2 GB/sec
相互結合網	Product		Intel® Omni-Path Architecture
	リンク速度		100 Gbps
	トポロジ		フルバイセクションバンド幅Fat-tree網





Intel® Omni-Path Architecture を用いた ^{CO JCAHPC} フルバイセクションバンド幅Fat-tree網



- コストはかかるがフルバイセクションバンド幅を維持
- ・システム全系使用時にも高い並列性能を実現
- 柔軟な運用:ジョブに対する計算ノード割り当ての自由度が高い







Oakforest-PACS の仕様(続き)

並列ファイルシステム	Туре	Lustre File System
	総容量	26.2 PB
	Product	DataDirect Networks SFA14KE
	総バンド幅	500 GB/sec
高速ファイ ルキャッ	Type	Burst Buffer, Infinite Memory Engine (by DDN)
シュシステ	総容量	940 TB (NVMe SSD, パリティを含む)
厶	Product	DataDirect Networks IME14K
	総バンド幅	1,560 GB/sec
総消費電力		4.2MW(冷却を含む)
総ラック数		102







Oakforest-PACS のソフトウェア

• OS:

- Red Hat Enterprise Linux (ログインノード)、
 CentOS および McKernel (計算ノード、切替可能)
- McKernel: 理研AICSで開発中のメニーコア向けOS
 - Linuxに比べ軽量、ユーザプログラムに与える影響なし
 - ポスト京コンピュータにも搭載される予定。

・コンパイラ

- GCC, Intel Compiler, XcalableMP
- XcalableMP:
 - 理研AICSと筑波大で共同開発中の並列プログラミング言語
 - CやFortranで記述されたコードに指示文を加えることで、性能の高い並列アプリケーションを簡易に開発することができる。

アプリケーション:

• OpenFOAM, ABINIT-MP, PHASE system, FrontFlow/blueなど、 オープンソースソフトウェア







計算ノードの写真





2Uサイズのシャーシ (富士通 PRIMERGY CX600 M1)に 8計算ノードを搭載

計算ノード (富士通 PRIMERGY CX1640 M1) Intel Xeon Phi 1ソケット、Intel Omni-Path Architecture card (HFI)搭載





運用予定

- ・スケジュール
 - 2016/10/1: 第1段階のシステム稼働(全系システムの5% 程度の規模)
 - 2016/12/1: 第2段階のシステム稼働(全系システム)
 - 2017/4:オープンな資源提供(HPCI資源を含む)
- 運用形態
 - ・通常運用:ハードウェアの分割ではなく、「CPU時間」を2 大学で按分することで柔軟な運用を可能に
 - 特別運用:限られた時間だけ、全系を1システムとして、超大規模な単一ジョブの実行(ex. Gordon Bell Challenge)
 - ・省電力運用:夏季など、状況に応じて、総電力にキャッピングをかける省電力運用







おわりに

- JCAHPC(最先端共同HPC基盤施設)
- ・筑波大学計算科学研究センターと東京大学情報基盤センターが設置
 - 計算科学・工学及びその推進のための計算機科学・工学の発展に資 するために連携して設置
- Oakforest-PACS:ピーク性能 25 PFLOPS
 - Intel Xeon Phi (Knights Landing) & Omni-Path Architecture
 - CPU時間を2大学で按分することで柔軟な運用を可能
 - 全系を1システムとして超大規模単一ジョブの実行も可能に
 - 全系システムの稼働は2016/12を予定
 - HPCI資源を含めオープンは資源提供は2017/4を予定
- JCAHPC: 最先端HPC研究に寄与する計算資源の提供を目指し、コミュニティに貢献していく予定



